

1 问题描述

2 解决方案¹

2.1 数据集的建立及模型的选择

2.2 探索模型单独工作的性能

1. 直接划分训练集和测试集
无可奉告
2. 交叉验证
3. 网格搜索

2.3 探索模型组合工作的性能

3 结果展示

3.1 模型单独工作的性能

- 3.1.1 直接划分训练集和测试集
- 3.1.2 交叉验证
- 3.1.3 网格搜索

3.2 模型组合工作的性能

- 3.2.1 **bagging**
- 3.2.2 **boosting**

¹本次作业的所有代码实现可参见附录 A.3

A 附录

A.1 Boston 数据集特征信息

表 1: Boston 数据集特征信息

编号	特征名	特征含义
1	CRIM	城镇人均犯罪率
2	ZN	住宅用地超过 25000 sq.ft. 的比例
3	INDUS	城镇非零售商用土地的比例
4	CHAS	查尔斯河空变量（如果边界是河流，则为 1；否则为 0）
5	NOX	一氧化氮浓度
6	RM	住宅平均房间数
7	AGE	1940 年之前建成的自用房屋比例
8	DIS	到波士顿五个中心区域的加权距离
9	RAD	辐射性公路的接近指数
10	TAX	每 10000 美元的全值财产税率
11	PTRATIO	城镇师生比例
12	B	$1000(B_k - 0.63)^2$, 其中 B_k 指代城镇中黑人的比例
13	LSTAT	人口中地位低下者的比例

A.2 PCA 算法在主成分个数为 3 时的进一步可视化

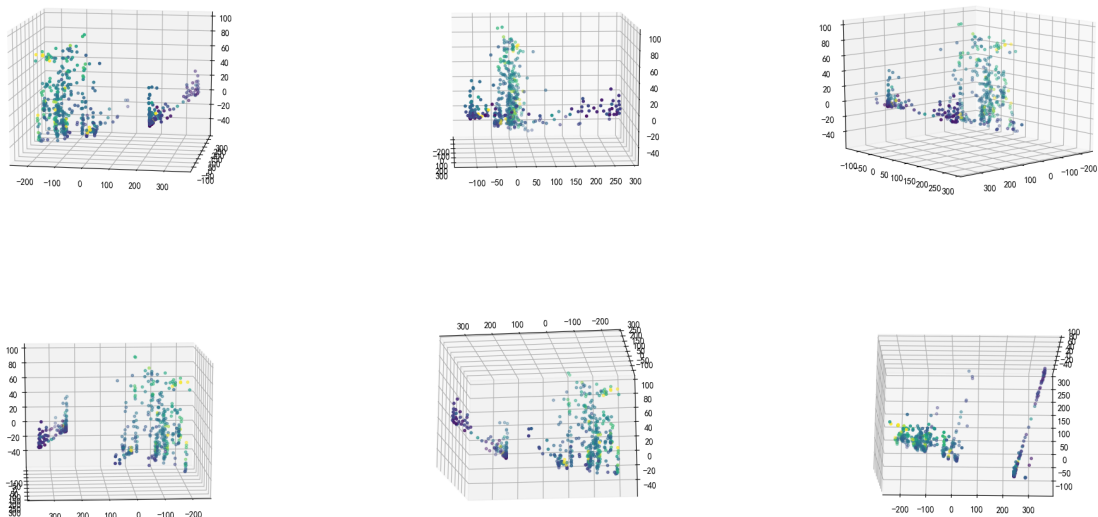


图 1: PCA 效果进一步可视化（主成分个数为 3 时）

A.3 主要代码

```

1 import matplotlib.pyplot as plt
2 plt.rcParams['font.sans-serif'] = ['SimHei']      # display Chinese
3 plt.rcParams['axes.unicode_minus'] = False        # display minus sign
4 from mpl_toolkits.mplot3d import Axes3D
5 from sklearn import datasets
6 from sklearn.decomposition import PCA
7
8 import os
9 import time
10
11 class bostonAnalyzer(object):
12     def __init__(self):
13         # load dataset
14         self.dataset = datasets.load_boston()
15         self.data = self.dataset.data
16         self.target = self.dataset.target
17
18         # for plot
19         self.var = []
20         self.t = []
21
22         if not os.path.exists('report/img'):
23             os.makedirs('report/img')
24
25     def run(self, vis=False):
26         '''
27         Run PCA for 1-13 dimensions
28         :param vis: True: plot dim 1-3; False: not plot
29         '''
30
31         self.var.clear()
32         self.t.clear()
33
34         with open('report/result.txt', 'w') as f:
35             for i in range(13):
36                 t = time.time()
37                 pca_op = PCA(n_components=i)
38                 pca_res = pca_op.fit_transform(self.data)
39                 t = time.time() - t
40
41                 self.var.append(pca_op.explained_variance_ratio_.sum()
42                                )
43                 self.t.append(t)
44
45                 # write log
46                 f.write('#####_Dimension_%d_#####\n' % i)
47                 f.write(str(pca_res.shape) + '\n')

```

```

47         f.write(str(pca_op.explained_variance_ratio_) + '\n')
48         f.write(str(pca_op.explained_variance_ratio_.sum()) +
49                 '\n')
50         f.write(str(t) + '\n\n')
51
52         # visualize for debug & plot
53         if vis:
54             if i == 1:
55                 # plot 1 dimension
56                 plt.scatter(pca_res[:, 0], pca_res[:, 0], s
57                             =14, c=self.target)
58                 plt.savefig('report/img/pca-%d' % i)
59                 plt.show()
60             elif i == 2:
61                 # plot 2 dimensions
62                 plt.scatter(pca_res[:, 0], pca_res[:, 1], s=8, c
63                             =self.target)
64                 plt.savefig('report/img/pca-%d' % i)
65                 plt.show()
66             elif i == 3:
67                 # plot 3 dimensions
68                 ax = plt.subplot(projection='3d')
69                 ax.scatter(pca_res[:, 0], pca_res[:, 1],
70                           pca_res[:, 2], s=8, c=self.target)
71                 plt.savefig('report/img/pca-%d' % i)
72                 plt.show()
73
74     def show(self):
75         '''
76         Show the basic info of Boston dataset & plot k-lines
77         '''
78
79         print(self.data.shape)
80         print(self.target.shape)
81         self._k_line_radio()
82         self._k_line_time()
83
84     def _k_line_radio(self):
85         '''
86         Plot k-line of variant radio
87         '''
88
89         x = list(range(len(self.var)))
90         plt.scatter(x, self.var, s=14, c='r')
91         plt.plot(x, self.var)
92         plt.xlabel('主成分个数')
93         plt.ylabel('降维后各特征方差比例之和')
94         plt.savefig('report/img/kline-radio')
95         plt.show()

```

```
91     def _k_line_time(self):
92         '''
93         Plot k-line of time
94         '''
95         x = list(range(len(self.t)))
96         plt.scatter(x, self.t, s=14, c='r')
97         plt.plot(x, self.t)
98         plt.xlabel('主成分个数')
99         plt.ylabel('PCA算法消耗时间(s)')
100        plt.savefig('report/img/kline-time')
101        plt.show()
102
103 if __name__ == '__main__':
104     bA = bostonAnalyzer()
105     bA.run(True)
106     bA.show()
```